

Experimental design and model selection: The example of exoplanet detection

Vijay Balasubramanian^{1,2}, Klaus Larjo¹ and Ravi Sheth^{1*}

¹*David Rittenhouse Laboratories, University of Pennsylvania,
Philadelphia, PA 19104, USA*

²*School of Natural Sciences, Institute for Advanced Study,
Princeton, NJ 08540, USA*

Abstract

We apply the Minimum Description Length model selection approach to the detection of extra-solar planets, and use this example to show how specification of the experimental design affects the prior distribution on the model parameter space and hence the posterior likelihood which, in turn, determines which model is regarded as most ‘correct’. Our analysis shows how conditioning on the experimental design can render a non-compact parameter space effectively compact, so that the MDL model selection problem becomes well-defined.

*vijay@physics.upenn.edu, klarjo@physics.upenn.edu, shethrk@physics.upenn.edu

1 Introduction

The Bayesian approach to parametric model selection requires the specification of a prior probability distribution over the parameter space. The Jeffreys’ prior, which is proportional to the square root of the determinant of the Fisher information computed in the parameter space, has been shown to be the uniform prior over all *distributions* indexed by the parameters in a parametric family [1]. Geometrically, its integral over a region of the parameter space computes a volume that essentially measures the fraction of statistically distinguishable probability distributions within that region [1]. In this interpretation, the Jeffreys prior distribution

$$\omega(\Theta) = \frac{\sqrt{\det J_{ij}(\Theta)}}{\int d^d\Theta \sqrt{\det J_{ij}(\Theta)}} d^d\Theta \quad (1)$$

where $\Theta = \{\theta_1, \dots, \theta_d\}$ simply measures the fractional volume of the small element $d^d\Theta$ relative to total volume of the parametric manifold $V = \int d^d\Theta \sqrt{\det J_{ij}(\Theta)}$. Here J_{ij} is the Fisher information on the parameter space $\Theta \in \mathbb{R}^d$ and $d^d\Theta$ is the standard Riemannian volume element on \mathbb{R}^d . The volume V also appears in the Minimum Description Length (MDL) approach to model selection [2, 3], conceptually because it effectively measures how many different distributions are describable by different parameter choices.

An important difficulty in applying the MDL approach to model selection occurs when the parameter space is noncompact and the volume V diverges. In this case, from the Bayesian perspective, a uniform prior on the parameter space does not exist, while from the MDL perspective the number of models that might be describable diverges, leading to problems with the definition of the description length. Of course the parameter space can be cut off by hand, but unless the choice of cut-off is well founded, it can lead to artifacts in the comparison of different model families [4, 5, 6]. Unfortunately in many practical problems the parameter space *is* noncompact and V diverges. For example, in astrophysics, the detection of exoplanets depends on a model of the light coming from the occluded star. This model will contain a non-compact direction representing the orbital period of the planet – see, e.g., [7]. For examples from psychophysics see, e.g., [4].

In this note we argue that merely specifying the experimental set-up – before the measurement of any actual data – influences the prior distribution on the parameter space. This occurs because, given the finite number of measurements in any experiment, many of the probability distributions indexed by a parametric manifold will be statistically indistinguishable. In cases where the parameter space is noncompact, the uniform prior conditioned on the experimental setup can thus become well-defined. In the geometric language of [1], the volume that measures the number of probability distributions in the parametric family that are statistically distinguishable given a *finite* number of measurements can be finite even if the parameter space is non-compact. In effect, specifying the experimental set-up can render the parameter space compact.

Our results illustrate how the choice of experimental set-up influences the measure on the parameter space of a model, thereby affecting which model is regarded as most ‘correct’. In

section 2 we briefly review the computation of posterior probabilities, and consider the effect of conditioning on the experimental set-up on the parameter space measure. In section 3 we apply these considerations to a physical problem: the analysis of light-curves of stars with orbiting planets. In this example we see that the volume of the parameter space is rendered effectively finite after the experimental set-up is specified.

2 The effect of experimental design on the parameter space measure

2.1 Review

Suppose one is interested in some physical phenomenon, and has made N relevant measurements: $Y = \{y_1, \dots, y_N\}$. Further suppose that there are two different parametric models, A and B , that aim to describe the phenomenon in question. The basic question to be answered is which of the two models is the better one, considering the experimental data Y . The probability-theoretic answer to this question is to compute the posterior probabilities $P(A|Y)$ and $P(B|Y)$, which we can write using the Bayes Rule as

$$P(A|Y) = \frac{P(A)}{P(Y)} \int \omega(\Theta) P(Y|\Theta), \quad (2)$$

where $\Theta = (\theta_1, \dots, \theta_d) \in \mathbb{R}^d$ is the vector of variables parametrising A , and $\omega(\Theta)$ is the volume form associated to the measure on the parameter space, which we will define shortly. A corresponding expression can also be written for $P(B|Y)$. Since we wish to compare $P(A|Y)$ and $P(B|Y)$, we can ignore the common factor $P(Y)$, and we will assume $P(A) = P(B)$ and drop this factor as well. Thus the only remaining ingredient to be defined is the volume form $\omega(\Theta)$; we simply quote the result from [1]: the volume form that gives equal weight to all statistically distinguishable distributions in the parametric family is

$$\omega(\Theta) = \frac{\sqrt{\det J_{ij}(\Theta)}}{\int d^d\Theta \sqrt{\det J_{ij}(\Theta)}} d^d\Theta, \quad (3)$$

where $J_{ij}(\Theta)$ is the *Fisher information matrix*, defined as the second derivative of the Kullback–Leibler distance $D(\Theta_p||\Theta_q)$:

$$J_{ij}(\Theta_p) = \partial_{\theta_i} \partial_{\theta_j} D(\Theta_p||\Theta_p + \Phi)|_{\Phi=0}, \quad (4)$$

$$D(\Theta_p||\Theta_q) = \int d\vec{x} \Theta_p(\vec{x}) \ln \frac{\Theta_p(\vec{x})}{\Theta_q(\vec{x})}. \quad (5)$$

where $d\vec{x}$ is the integration measure over the sample space $\{\vec{x}\}$, and $\Theta_p(\vec{x})$ is the distribution function associated to the values of the parameters $(\theta_1^p, \dots, \theta_d^p)$. Now we have defined everything needed to compute the posterior probabilities, and we illustrate the formalism by applying it to the analysis of light-curves.

Using this, we can compute the Fisher information matrix by computing the Kullback–Leibler distance between two nearby points and Taylor expanding:

$$\begin{aligned}
D(\Theta_0||\Theta_q) &= \int d^N \vec{y} \Theta_0(\vec{y}) \ln \frac{\Theta_0(\vec{y})}{\Theta_q(\vec{y})} \\
&\approx - \int d^N \vec{y} \Theta_0(\vec{y}) \ln \frac{\Theta_0(\vec{y}) + \partial_{\theta_i} \Theta_0(\vec{y}) \Delta \theta_i + \partial_{\theta_i} \partial_{\theta_j} \Theta_0(\vec{y}) \Delta \theta_i \Delta \theta_j}{\Theta_0(\vec{y})} \\
&\approx - \int d^N \vec{y} \left(\partial_{\theta_i} \Theta_0(\vec{y}) \Delta \theta_i + \partial_{\theta_i} \partial_{\theta_j} \Theta_0(\vec{y}) \Delta \theta_i \Delta \theta_j - \frac{1}{2} \frac{(\partial_{\theta_i} \Theta_0(\vec{y}))(\partial_{\theta_j} \Theta_0(\vec{y}))}{\Theta_0(\vec{y})} \Delta \theta_i \Delta \theta_j \right) \\
&= \frac{1}{2} \int d^N \vec{y} \underbrace{\frac{(\partial_{\theta_i} \Theta_0(\vec{y}))(\partial_{\theta_j} \Theta_0(\vec{y}))}{\Theta_0(\vec{y})}}_{\equiv J_{ij}(\Theta_0)} \Delta \theta_i \Delta \theta_j. \tag{6}
\end{aligned}$$

On the third line, the terms linear in Θ_0 vanish, as exchanging the order of integration and derivation, the integral of Θ_0 will yield a constant 1, which then differentiates to zero.

2.2 Effect of the experimental set-up

The measure (3) is independent of the experimental data Y and is constructed under the assumption that the entire sample space can be measured by the observer. However, in real experiments, instrumental and design limitations only allow observation of some subset M of the sample space. Thus an observation either results in no detected outcome, or in a measurement $y_i \in M$. Thus the effective predicted distribution of measured outcomes is not the $\Theta(\vec{y})$, but rather

$$\Theta(\vec{y}) = \begin{cases} \Theta(\vec{y}), & \text{for } \vec{y} \in M, \\ \Theta^{\text{Out}}, & \text{no measured outcome,} \end{cases} \tag{7}$$

where $\Theta^{\text{Out}} \equiv \int_{\vec{y} \notin M} d\vec{y} \Theta(\vec{y})$. We will argue that if the models in the asymptotic regions of a noncompact parameter space differ in their predictions mostly outside the observable region M , the Fisher information for the effective distributions (7) can decay sufficiently quickly to render the volume $V = \int d^d \Theta \sqrt{\det J_{ij}(\Theta)}$ finite. In this section we will give one set of sufficient conditions for this to happen and in Sec. 3 we will give a detailed example.

Consider a model, specified by parameters $\vec{\theta} = (\theta_1, \dots, \theta_d) \in \mathbb{R}^d$, and a distribution $\Theta_{\vec{\theta}}(\vec{x})$, with $\vec{y} \in \mathbb{R}^n$. We will slightly simplify notation simply referring to the distribution as $\Theta(\vec{y})$ and understanding the implicit parameter dependence. Let us use spherical coordinates in the parameter space \mathbb{R}^d with ρ being the radial coordinate, i.e. $(\theta_1, \dots, \theta_d) \rightarrow (\rho, \varphi_1, \dots, \varphi_{d-1})$. Also consider an experimental set-up that can only make measurements inside some compact region $M \subset \mathbb{R}^n$. Thus, the probability of no measurement being registered by this experiment is $\Theta^{\text{Out}} \equiv \int_{\vec{y} \notin M} d\vec{y} \Theta(\vec{y})$.

Our first assumption is a smoothness condition, so that inside the region M the distribution does not fluctuate too much as one approaches the asymptotics of parameter space:

$$\left| \partial_i \Theta(\vec{y})|_{\vec{y} \in M} \right| \leq \delta(\rho), \quad \text{for large } \rho, \quad i = 1, \dots, d, \tag{8}$$

where $\delta(\rho)$ goes to zero as ρ goes to infinity; we will later specify the exact scaling needed. Intuitively, this condition says that as the parameter $\rho \rightarrow \infty$, the models do not differ too much inside the observable part of the sample space M . This allows us to estimate

$$|\partial_i \Theta^{\text{Out}}| = \left| \partial_i \left(1 - \int_{y \in M} d\vec{y} \Theta(\vec{y}) \right) \right| \leq \text{Vol}(M) \delta, \quad (9)$$

where $\text{Vol}(M)$ denotes the volume of the compact region M .

Secondly we assume that inside M , the distributions $\Theta(\vec{y})$ do not decay too quickly as $\rho \rightarrow \infty$. Intuitively, since any experiment will only measure a finite amount of data (say N points), if the probability of a single measurement lying inside M is significantly less than $1/N$, then the experimental set-up will not detect anything. Thus we will require

$$\Theta(\vec{y})|_{\vec{y} \in M} > \epsilon(\rho), \quad \text{for large } \rho, \quad (10)$$

where again we will later specify the scaling of $\epsilon(\rho)$ with ρ .¹

Using these assumptions, we can establish an upper bound for the Fisher information (6):

$$\begin{aligned} |J_{ij}| &\leq \left| \int_{\vec{y} \in M} \frac{\partial_i \Theta(\vec{y}) \partial_j \Theta(\vec{y})}{\Theta(\vec{y})} \right| + \left| \frac{\partial_i \Theta^{\text{Out}} \partial_j \Theta^{\text{Out}}}{\Theta^{\text{Out}}} \right| \\ &< \delta^2 \left| \int_{\vec{y} \in M} \frac{1}{\Theta(\vec{y})} \right| + \text{Vol}(M)^2 \delta^2 \leq \text{Vol}(M) \frac{\delta^2}{\epsilon} + \text{Vol}(M)^2 \delta^2 \sim \text{Vol}(M) \frac{\delta^2}{\epsilon}. \end{aligned} \quad (11)$$

Thus the determinant of the Fisher information scales as

$$\sqrt{\text{Det } J_{ij}} \sim \left(\frac{\delta^2}{\epsilon} \right)^{\frac{d}{2}}, \quad (12)$$

and for the integral V to be finite one must have suppression stronger than $\sqrt{\text{Det } J_{ij}} \sim \rho^{-d}$. Thus the integral converges if δ is suppressed more strongly than

$$\delta(\rho) < \frac{\sqrt{\epsilon(\rho)}}{\rho}. \quad (13)$$

From the experimental set-up one can estimate how $\epsilon(\rho)$ scales with ρ , which then determines how $\delta(\rho)$ needs to scale for the integral to converge. This is thus a sufficient condition for rendering the parameter space effectively finite.

It is worth stressing that, following the above analysis, any method of deciding the validity of a model is impacted by the choice of the experiment in a completely computable way, and this should be taken into account when designing experiments.

¹This condition can be relaxed by recognizing that if $\Theta(\vec{y})|_{\vec{y} \in M}$ decays too quickly as $\rho \rightarrow \infty$, then the models in the asymptotic region of the parameter space make no measurable predictions for experiments designed with a finite number of measurements. The example in the Sec. 3 will illustrate such a scenario.

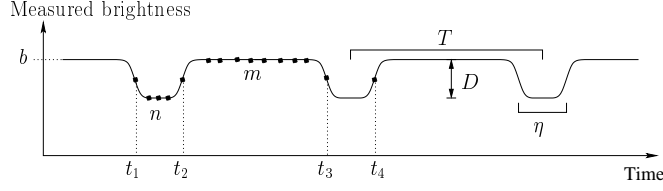


Figure 1: An example of a light-curve.

3 The probability of exo-planet detection

3.1 Model for exo-planets

Consider a star orbited by a planet so that the planet periodically passes between the star and Earth. The light output (light-curve) of such a star is a constant line, with a small periodic dip when the planet is eclipsing part of the star. One model for such a light-curve was proposed in [7] as

$$y(T, D, \eta, \tau, b; t) = b - \frac{D}{2} \left[\tanh c\left(\tilde{t} + \frac{1}{2}\right) - \tanh c\left(\tilde{t} - \frac{1}{2}\right) \right], \quad (14)$$

where

$$\tilde{t} = \frac{T \sin \frac{\pi(t-\tau)}{T}}{\pi\eta}. \quad (15)$$

An example light-curve is shown in figure 1; T is the period of the planet; η is the duration of the transit, i.e. how long the planet eclipses the star; D is the depth of the dip in the curve; b is the total observed brightness of the star; and τ is a phase parameter specifying when the planets transit occurs. Finally, c is a constant parameter specifying the sharpness of the edges of the light-curve, expected to be fairly large as the transition between transit/no-transit is relatively quick. The assumption $c \gg 1$ greatly simplifies our analysis, and is not physically very restrictive.

The parameter space for this model is clearly non-compact as T can range to infinity. However, we will argue that the space is effectively rendered compact after the experimental set-up is specified. To be precise, the parameter space is²:

$$T \in [0, \infty), \quad D \in [0, b], \quad \tau \in [0, T], \quad \eta \in [0, \delta T], \quad b \in [0, b_{max}], \quad (16)$$

where δ is a small number that we will estimate, and the maximal brightness b_{max} is naturally given by the brightness of Sirius, the brightest star visible from Earth. Assuming a circular orbit as in Figure 2, the ratio of the transit time to the period of the planet is given by

$$\frac{\eta}{T} \approx \frac{2r/v_{\text{planet}}}{2\pi R/v_{\text{planet}}} = \frac{1}{\pi} \frac{r}{R}.$$

²Note that we consider c to be a constant, not a parameter.

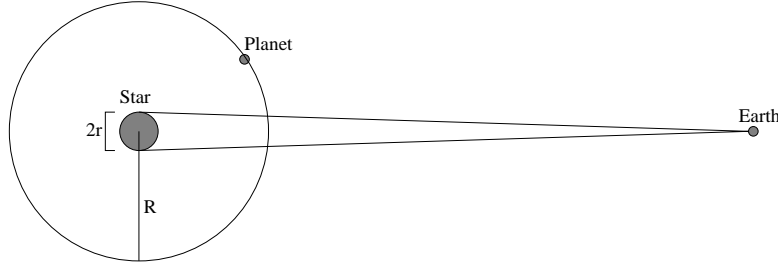


Figure 2: The basic set-up: an extra-solar planet orbiting a star of radius r with an average distance R .

For the currently known transiting exo-planets this ratio is around ~ 0.1 [8], although for a typical system one expects it to be smaller as large planets orbiting close to the star are easier to observe, which favors largest values of the ratio. For an elliptical orbit, the answer will differ by an $\mathcal{O}(1)$ factor, but will have the same dependence on r/R . Thus, η will always be a small fraction of T .

Now we can write down the probability density for measuring values $\vec{y} = (y_1, \dots, y_N)$ for the light-curve at times (t_1, \dots, t_N) with the light-curve specified by parameters $(\theta_1^0, \theta_2^0, \theta_3^0, \theta_4^0, \theta_5^0) = (T, D, \eta, \tau, b)$ as

$$\Theta_0(\vec{y}) = \prod_{k=1}^N \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{(y_k - y_0(\theta_i^0; t_k))^2}{2\sigma_k^2}} = (2\pi\sigma)^{-\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \sum_{k=1}^N (y_k - y_0(\theta_i^0; t_k))^2}, \quad (17)$$

where we have assumed that the uncertainty in each measurement is Gaussian, and further we have chosen the standard deviation to be equal for all measurements for simplicity. Using (17) in the formula (6), we see that the the integrals in the Fisher information are Gaussian in y_k ; thus we can compute them analytically to get

$$J_{ij} = \frac{1}{\sigma^2} \sum_{k=1}^N \partial_{\theta_i} y(\theta; t_k) \partial_{\theta_j} y(\theta; t_k). \quad (18)$$

This is our key formula, and we shall spend the next subsection analysing its properties.

3.2 Finiteness of light-curve parameter space

We now wish to apply the general arguments of section 2 to the exo-planet system. Consider an experimental set-up that can barely measure two periods, and then consider shortening the experiment slightly so that only one dip is detected; this is depicted in figure 1. To be precise, the shorter set-up measures the beginning and end of a transit at t_1 and t_2 , n points in between, and m points after the transit. The longer set-up makes measurements at the same times, and additionally at times t_3 and t_4 , detecting the second transit. In the next subsections we will show that $J_{\text{short}} \lll J_{\text{long}}$, indicating that detecting the second dip is of

fundamental importance to experimental design; without the second dip the experimental set-up can't differentiate models with large enough T . This renders the parameter space effectively finite, as an experiment can not differentiate between models that have period T larger than the duration of the experiment.

3.2.1 Effect of measuring a second transit on $\det(J)$

In this subsection we will give an estimate for the magnitude of the determinant of the Fisher information, and show how it is affected by the inclusion of the second transit in the data. In subsequent subsections we will exactly compute the determinant for a few specific experimental set-ups.

From (18) and the definition of a determinant, we see that in each term of the determinant each parameter θ_i appears exactly twice in the derivatives, i.e. each term is of the form

$$J_{i_1 j_1} J_{i_2 j_2} J_{i_3 j_3} J_{i_4 j_4} J_{i_5 j_5} \sim \frac{1}{\sigma^{10}} \partial_{Ty} \partial_{Ty} \partial_{Dy} \partial_{Dy} \partial_{\eta y} \partial_{\eta y} \partial_{\tau y} \partial_{\tau y} \partial_{by} \partial_{by}. \quad (19)$$

As a rough estimate of the size the determinant, we investigate how large terms of this type can be. The derivatives are

$$\partial_{Ty}(\theta, t) = \frac{cD}{2} \frac{f(\tilde{t})}{\pi\eta} \left(\sin \frac{\pi(t-\tau)}{T} - \frac{\pi(t-\tau)}{T} \cos \frac{\pi(t-\tau)}{T} \right), \quad (20)$$

$$\partial_{Dy}(\theta, t) = \frac{1}{2} \left(\tanh c(\tilde{t} - \frac{1}{2}) - \tanh c(\tilde{t} + \frac{1}{2}) \right), \quad (21)$$

$$\partial_{\tau y}(\theta, t) = -\frac{cD}{2} \frac{f(\tilde{t})}{\eta} \cos \frac{\pi(t-\tau)}{T}, \quad (22)$$

$$\partial_{\eta y}(\theta, t) = -\frac{cD}{2} \frac{\tilde{t} f(\tilde{t})}{\eta}, \quad \partial_{by}(\theta, t) = 1, \quad (23)$$

$$\text{with} \quad f(\tilde{t}) \equiv \tanh^2 c(\tilde{t} + \frac{1}{2}) - \tanh^2 c(\tilde{t} - \frac{1}{2}). \quad (24)$$

From (24) we see that $f(\tilde{t}) \neq 0$ only when $\tilde{t} \approx \pm \frac{1}{2}$, again assuming large c . This tells us that the measurements that contribute most to the Fisher information are the ones on the edges of the dips³, i.e. at times t_1, t_2, t_3 and t_4 in figure 1. We write the condition $|\tilde{t}| \approx \frac{1}{2}$ as

$$\left| \sin \frac{\pi(t-\tau)}{T} \right| = \frac{\pi}{2} \frac{\eta}{T}, \quad (25)$$

and note that the ratio of transit time to period is very small, $\frac{\eta}{T} \ll 1$. This gives us the solutions

$$\frac{t-\tau}{T} \approx n \pm \frac{\eta}{2T}, \quad (26)$$

where n is an integer indexing the number of the dip, with $n = 0$ denoting the solitary dip if only one is present in the data.

³This statement is somewhat subtle, and we will discuss this matter in more detail in section 3.2.3; for our current purposes it is sufficiently accurate.

We wish to estimate the ratio of the determinants of the Fisher information by an order of magnitude estimate

$$\frac{J_{\text{short}}}{J_{\text{long}}} \sim \frac{J_{i_1 j_1}^s J_{i_2 j_2}^s J_{i_3 j_3}^s J_{i_4 j_4}^s J_{i_5 j_5}^s |_{\text{max}}}{J_{i_1 j_1}^l J_{i_2 j_2}^l J_{i_3 j_3}^l J_{i_4 j_4}^l J_{i_5 j_5}^l |_{\text{max}}}, \quad (27)$$

where both the numerator and the denominator are of the form (19), and according to the argument above the maximal contributions come from the edge measurements. From (21-23) we see that the derivatives with respect to D, η, τ and b are all periodic at the edges: $|\partial_{\theta_i} y(\theta, t_1)| = \dots = |\partial_{\theta_i} y(\theta, t_4)|$ for $\theta_i \neq T$, and thus will cancel in the ratio (27).

It is crucial that $\partial_T y$, however, is not periodic due to the second term in (20). At the first dip, $t_1, t_2 = \pm \frac{\eta}{2T}$, we expand (20) to find

$$\partial_T y(t_1) \approx \partial_T y(t_2) \approx \frac{\pi^2 c D}{48} \frac{\eta^2}{T^3}, \quad (28)$$

while at the second dip, $t_3, t_4 = 1 \pm \frac{\eta}{2T}$, the contribution is

$$|\partial_T y(t_3)| \approx |\partial_T y(t_4)| \approx \frac{cD}{2\eta}, \quad (29)$$

ignoring signs that are irrelevant for this estimate. Thus we see that the Fisher information increases strongly as the second dip is included:

$$\frac{J_{\text{short}}}{J_{\text{long}}} \sim \frac{(\partial_T y(t_{1,2}))^2}{(\partial_T y(t_{1,2}))^2 + (\partial_T y(t_{1,2}) \partial_T y(t_{3,4})) + (\partial_T y(t_{3,4}))^2} \sim \left(\frac{\eta}{T}\right)^6 \lll 1, \quad (30)$$

where we ignored order one coefficients. This is an explicit example of how our arguments from section 2 work for a realistic model: when an experimental set-up does not have the capability to detect two dips, it becomes impossible to determine the period, and consequently the Fisher information is very small (or vanishing) compared to an experiment that is able to detect two dips and determine the period more accurately. For any given experiment of finite duration Δt , the Fisher Information will decline with T when $T \gg \Delta t$ effectively rendering the parameter space compact.

3.2.2 The tail $T \rightarrow \infty$

To verify our claim that the parameter space is really rendered compact we need to show that $\det J \rightarrow 0$ strongly enough as T is taken to infinity. It is easy enough to find the T -scaling of the derivatives (20-23); $\partial_T y$ scales as T^{-3} , while the others stay finite in the large T limit. Thus, as seen from (19), the determinant will scale as

$$\sqrt{\det J} \sim \partial_T y \sim \frac{1}{T^3}, \quad (31)$$

which shows that the parameter space measure vanishes fast enough for large T to render the parameter space volume finite.

3.2.3 Explicit computation of $\text{Det}(J_{ij})$ for specific experimental set-ups

While the order of magnitude estimate of the previous subsection offers an intuitive reason as to why the Fisher information decreases sharply when the number of peaks detected falls below two, it is still instructive to explicitly compute the determinant in a few experimental set-ups.

Detecting two dips: Let us first consider the case J_{long} from section 3.2, i.e. measurements at times indicated in figure 1. Using the derivatives (20-23) one can write down the Fisher information matrix (18) as

$$J_{ij}^{\text{long}} = \begin{pmatrix} 2(T_1^2 + T_3^2) & -T_1 & 2T_1X & -4T_3X & 2T_1 \\ -T_1 & 1+n & -2X & 0 & -(2+n) \\ 2T_1X & -2X & 4X^2 & 0 & 4X \\ -4T_3X & 0 & 0 & 16X^2 & 0 \\ 2T_1 & -(2+n) & 4X & 0 & 4+n+m \end{pmatrix}, \quad (32)$$

where for brevity we defined

$$T_1 \equiv \partial_T y(t_1) = \partial_T y(t_2) = \frac{cD\pi^2}{48} \frac{\eta^2}{T^3}, \quad T_3 \equiv \partial_T y(t_3) = -\partial_T y(t_4) = \frac{cD}{2\eta}, \quad (33)$$

$$X \equiv -\frac{cD}{4\eta} = \partial_\eta y(t_{1,2,3,4}) = -\frac{\partial_\tau y(t_{1,3})}{2} = \frac{\partial_\tau y(t_{2,4})}{2}. \quad (34)$$

In computing this matrix we used that $f(\tilde{t}) = 0$ for $\tilde{t} \neq \pm\frac{1}{2}$, which is true up to corrections of order e^{-c} , as seen from (24); for this reason one does not need to specify the exact times of the n measurements during the dip, or the m measurements outside the dip, as up to e^{-c} corrections they all contribute equally. The determinant of the Fisher information is simple,

$$\text{Det}(J_{ij}^{\text{long}}) = 64nmX^4(T_1^2 + T_3^2) \approx 64nmX^4T_3^2. \quad (35)$$

This result explains the subtlety referred to earlier: although measurements at the edges contribute the most to the Fisher information, if one only has measurements at the edges ($n = m = 0$) the Fisher information actually vanishes. Physically this is easy to interpret, as only measuring the edges t_1, \dots, t_4 will yield four points lying on a line, and thus they cannot be used to determine any information about the curve; other data points are needed to ‘anchor’ the data.

Detecting only one dip: Similarly one can compute the Fisher information in the ‘short’ experimental set-up, where measurements are made at the same times as before, except not at t_3 and t_4 . This yields

$$J_{ij}^{\text{short}} = \begin{pmatrix} 2T_1^2 & -T_1 & 2T_1X & 0 & 2T_1 \\ -T_1 & \frac{1}{2} + n & -X & 0 & -(1+n) \\ 2T_1X & -X & 2X^2 & 0 & 2X \\ 0 & 0 & 0 & 8X^2 & 0 \\ 2T_1 & -(1+n) & 2X & 0 & 2+n+m \end{pmatrix}, \quad (36)$$

and perhaps surprisingly the determinant vanishes: $\text{Det}(J_{ij}^{\text{short}}) = 0$, up to tiny e^{-c} corrections. This indicates that the estimate in section 3.2.1 was an overestimate⁴: terms in the determinant of J^{short} are of the magnitude estimated, but the determinant is arranged in such a way that the terms cancel to a high accuracy, and the compactness of the parameter space is strengthened.

4 Discussion

Our analysis has shown how the specification of an experimental design affects the measure on model parameter spaces in MDL model selection (or equivalently the prior probability distribution on parameters in the Bayesian approach). Interestingly, the finite number of measurements within a bounded sample space in any practical experiment can effectively render a non-compact parameter space compact thereby leading to a well-defined prior distribution (3). Our analysis could be turned around to design experiments to discriminate well between models in some chosen region of the parameter space by ensuring that the Fisher information (18) is large in the desired region. It would also be useful to determine general conditions under which experimental design effectively makes model parameter spaces compact, perhaps following the arguments of Sec. 2.

Acknowledgments: This paper was written in honor of Jorma Rissanen’s 75th birthday and his many seminal achievements in statistics and information theory. VB and KL were partially supported by the DOE under grant DE-FG02-95ER40893, and KL was also partly supported by a fellowship from the Academy of Finland. VB was also partly supported as the Helen and Martin Chooljian member at the Institute for Advanced Study.

References

- [1] V. Balasubramanian, “Statistical Inference, Occam’s Razor and Statistical Mechanics on The Space of Probability Distributions,” [arXiv:cond-mat/9601030], V. Balasubramanian, “A Geometric Formulation of Occam’s Razor for Inference of Parametric Distributions,” [arXiv:adap-org/9601001].
- [2] J. Rissanen, “Modeling by shortest data description”, *Automatica*, 14:1080-1100, 1978.
- [3] J. Rissanen, “Fisher information and stochastic complexity”, *IEEE Trans. Inform.Theory*, 42:40-47, 1996.
- [4] I.J. Myung, V. Balasubramanian and M.A. Pitt, “Counting Probability Distributions: Differential Geometry and Model Selection”, *Proceedings of the National Academy of Science*, 97(21) 11170–11175, 2000.

⁴As the estimate illustrates an intuitive reason why the appearance of the second peak is so important, we decided to include it.

- [5] F. Liang and A. R. Barron, “Exact minimax strategies for predictive density estimation, data compression, and model selection”, IEEE Transactions on Information Theory 50, 2708-2726, 2004.
- [6] Chapter 11 of P.D. Grünwald, *The Minimum Description Length Principle*, MIT Press, June 2007.
- [7] P. Protopapas, R. Jimenez and C. Alcock, “Fast identification of transits from light-curves,” Mon. Not. Roy. Astron. Soc. **362**, 460 (2005) [arXiv:astro-ph/0502301].
- [8] <http://obswww.unige.ch/~pont/simpleTABLE.dat>